

1. A perceptron is linear, but XOR is not linearly separable.
MLPs add hidden layers and nonlinearity, enabling multiple decision boundaries.
2. Even though the composition of linear function is still linear:
Gradients shrink because of multiple repeated multiplication of derivatives < 1 like sigmoid or tanh.
ReLU doesn't saturate for positive inputs & has derivative 1.
3. Necessary because self attention looks sequence order.
Sinusoidal PE is fixed & extrapolates to longer sequences.
Absolute PE is flexible but limited to trained sequence.
ROPE encodes relative position in attention space
4. Query - what a token is searching for
Key - what token represents.
Value - info returned after attention.

Scaling by $\sqrt{d_k}$ prevents large dot products that cause softmax saturation.
- Because tokens strongly attend to themselves.
5. Multiple heads allow the model to attend to different features & relationships simultaneously.

d_{model} - total embeddings dimension
h - no. of attention heads

d-head = $\frac{d_{model}}{h}$

- 6- Greedy decoding selects highest probability token at each step
 Beam keeps multiple candidate sequences.
 Eg - Greedy \rightarrow I am fine
 Beam \rightarrow I am finally feeling much better

$$1. \quad d_{head} = \frac{d_{model}}{h} = \frac{768}{12} = 64$$

Q, K, V has a weight matrix of size 768×768
 parameters per matrix = $768^2 = 589,824$
 Total for $Q, K, V = 3 \times 589,824 =$

$$2. \quad \text{Softmax} = \frac{e^{x_i}}{\sum e^{x_i}} \rightarrow e^2 + e^1 + e^0 = 11.11$$

$$\frac{e^2}{11.11} = 0.665 \quad \frac{e^1}{11.11} = 0.245 \quad \frac{e^0}{11.11} = 0.090$$

$$[0.665, 0.245, 0.090]$$

3. Numpy -

Import numpy as np

def attention (Q, K, V):

$dk = Q.shape[-1]$

Scores = $Q @ K.T / (np.sqrt(dk))$

weights = $np.exp(scores) / np.sum(np.exp(scores))$

output = weights @ V

return output, weights

axis=1, keepdim
= True

4. Numpy Version - Masked Attention.

import numpy as np

def masked_attention(Q, K, V):

dk = Q.shape[-1]

scores = Q @ K.T / np.sqrt(dk)

seq_len = scores.shape[0]

mask = np.triu(np.ones((seq_len, seq_len)), k=1)

scores = np.where(mask == 1, -np.inf, scores)

weights = np.exp(scores) / np.sum(np.exp(scores), axis=1)

~~output = weights @ V~~ keepdim = True)

return output, weights.